

User-intuitive Explanations for Increasing the Transparency of Autonomous Agents

Francesca Stramandinoli^{1*}, Brett Israelsen¹, Peggy Wu¹, Kishore K. Reddy¹, Franklin R. Tanner², Laura Strater²

¹Raytheon Technologies Research Center

²Raytheon Intelligence & Space

*** Correspondence:**

Corresponding Author

francesca.stramandinoli@rtx.com

Keywords: Explainable and Trusted AI, Black-box Explanation, Interpretability Methods, Model-agnostic Explainability, Interpretable Features, Concept-based Explanation, Neuro-symbolic AI.

Abstract

In high tempo, high stakes military applications of Artificial Intelligence (AI) and Machine Learning (ML), operators need to rapidly understand the strengths and limitations of their AI/ML aides, so that the combined human + AI team can provide decision superiority. As AI agents become more sophisticated, operators need transparency in these advanced systems much like the ability to build mental models of their human collaborators to ultimately judge their appropriate use.

A key challenge in making AI explainable to humans who are not computer or data scientists, and who do not have the time or tools to understand the inner workings of the AI system, is to map hidden layers of machine reasoning to semantics that are human interpretable to gain more insight into AI recommendations so that operators may either have more confidence in AI results or know when to override it. To address this challenge of aligning semantics of the machine generated explanation with human interpretation, we first describe the creation of a surrogate white-box approach as a stand-in, and subsequently describe the concept of a semantic alignment method. We describe a use case of Automatic Target Recognition (ATR) and illustrate why current explainer approaches are inadequate in achieving machine transparency and discuss research needs.

30 1 Introduction

31 The NORAD and U.S. Northern Command (USNORTHCOM) collect and coordinate a worldwide
32 system of sensors to monitor for aerospace and maritime threats to the United States and Canada. The
33 network of satellites, ground-based radar, airborne radars, and other assets provide a tremendous
34 amount of data. Making sense of this data requires advanced capabilities such as Automatic Target
35 Recognition (ATR) powered by AI/ML. However, autonomous systems alone do not provide the value
36 of human experience and intuition. Further, our adversaries have been monitoring our methods to deter,
37 compete, and conduct war over the past decades. To ensure information dominance and decision
38 superiority, there is a need to combine the ingenuity and agility of human expertise with the pattern
39 recognition and scalable processing power of AI/ML. Human-Autonomy Teaming (HAT) presents
40 unique challenges. Just like human teams, operators need the ability to continuously assess and reassess
41 another team member's performance within different contexts, whether the other team member is
42 human or digital. Current AI/ML methods lack the inherent mechanisms that human teams have in
43 order to represent their levels of competence and trustworthiness. Worse, they lack the ability to be
44 transparent to a human operator. We adopt the National Academies of Sciences (NAS)'s definition of
45 transparency [1] as "*the understandability and predictability of the system*" (Endsley, Bolte, and Jones,
46 2003, p. 146 [2]), including the AI system's "*abilities to afford an operator's comprehension about an*
47 *intelligent agent's intent, performance, future plans, and reasoning process*" (Chen et al., 2014a, p. 2
48 [3]).

49 Intelligent systems based on opaque AI/ML methods might hide potential issues inherited by biased
50 data or lead to adopting decisions that we do not fully understand, or even worse, that violate ethical
51 principles. Indeed, very often for achieving high performance in prediction, recommendation, and
52 decision-making support, the adoption of complex models that hide the logic of their internal processes
53 is required. Such models are often referred to as black box. A large amount of work that aimed at
54 enabling end users to better understand, trust, and effectively manage artificially intelligent systems
55 has been conducted under the Defense Advanced Research Projects Agency (DARPA) eXplainable
56 Artificial Intelligence (XAI) program [4]. The XAI program envisioned three general approaches for
57 improving explainability while maintaining a high level of learning performance: i) deep explanation
58 approaches to learn features or representations that can make the inner workings of a deep learning
59 model more transparent; ii) interpretable model approaches to develop models that are inherently more
60 explainable and more introspectable for machine learning experts; and iii) model induction approaches
61 that experiment with any given ML model - such as a black box - to infer an approximate explainable
62 model.

63 This work reviews existing methods for explaining black box models with a focus on techniques that
64 can generate user-intuitive explanations. We describe the work done towards a prototype
65 implementation for generating explanations for AI/ML models that operate on a single modality of
66 data and discuss why current explainer approaches are inadequate in achieving machine transparency
67 and discuss research needs. The use case driving this investigation is Automatic Target Recognition
68 (ATR) where a human analyst is working with the aid of an AI/ML target recognition algorithm.

69 This paper is outlined as follows. The rest of this section is dedicated to background and foundational
70 concepts, and it describes the particular use case we have had in mind while performing this work.
71 Section 2 is focused on an overview of interpretability methods in use. Section 3 describes the current
72 prototype implementation for explanation generation and research challenges that need to be addressed
73 in pursuit of creating detailed, effective, explanations of black box machine learning models. Section

74 4 outlines metrics that have been proposed for explainability. Finally, Section 5 contains conclusions,
75 discussion, and possible future research directions.

76 1.1 Background and Foundational Concepts

77 In this section we introduce the definitions of transparency, explanations, and interpretability adopted
78 in this work. As noted above, we adopt the National Academies of Sciences (NAS)’s definition of
79 transparency [1] as “*the understandability and predictability of the system*” (Endsley, Bolte, and Jones,
80 2003, p. 146 [2]), including the AI system’s “*abilities to afford an operator’s comprehension about an*
81 *intelligent agent’s intent, performance, future plans, and reasoning process*” (Chen et al., 2014a, p. 2
82 [3]). In other words, a model is transparent if a human can maintain an accurate mental model of how
83 the AI/ML system works. One way to improve model transparency is to provide valid and
84 comprehensible explanations to the human. Following the definition in [5] “*an explanation is an*
85 *“interface” between humans and a decision maker that is at the same time both an accurate proxy of*
86 *the decision maker and comprehensible to humans*”. In order to create explanations, interpretability
87 methods should be leveraged. We adopt Doshi-Velez and Kim’s definition of interpretability as “*the*
88 *ability to explain or to present in understandable terms to a human*” and extend upon it to add “*without*
89 *any additional machine processing*” [6]. In other words, an interpretable model is inherently self-
90 explanatory by an operator knowledgeable of the subject. We view interpretability as a continuum that
91 examines the extent to which a process or model can be understood by a human without any assistance
92 (if a user viewed the model alone could they make sense of it? This is discussed further in [37]). An
93 example of an interpretable system is a rule-based system where heuristics used by the model are
94 documented and made available to the human interpreter. This understanding of models is critical for
95 alignment with DoD’s Ethical Principles of Traceable and Reliable Artificial Intelligence [7].
96 Interpretability may not be required in all systems, especially those that can complete their tasks
97 without human intervention [6] (e.g., from toaster to GPS systems). However, creating interpretable
98 systems is necessary in high stakes safety critical human-in-the-loop systems where the ability to
99 understand the machine’s reasoning can identify misalignments in objectives and in turn improve a
100 user’s ability to make appropriate trust-based decisions.

101 The design of an interpretable system requires taking several factors into account. An important aspect
102 to consider is the user who is interested in receiving the explanations and why he/she needs an
103 explanation. Identifying the user group (e.g., decision maker, person affected by the decision, designer
104 of the system) helps determining the level of details and type of information to include in the
105 explanation. Another important aspect to consider are the questions that an interpretable system should
106 answer. Different types of AI / ML systems can answer different types of questions (e.g., why, why
107 not, what, what if, how, how confident, etc.). For example, an ML system trained to solve a
108 classification problem can answer why the model predicts a specific outcome, while a model trained
109 to solve a decision-making task can explain what caused a specific event, showing causal relationships.
110 Moreover, an explanation needs to be contextualized according to when the user needs it (e.g., during
111 onboarding, regular interactions, system errors / malfunctioning, etc.). In our current work we focus
112 on explainable systems for decision makers that need to have answers to why questions during regular
113 interactions with the system.

114 Interpretability methods can be classified along several dimensions [8]. An important distinction can
115 be made based on the type of algorithm it can be applied to. Model-specific approaches are specifically
116 designed with a model type in mind and therefore they generally require detailed information about the
117 model; perhaps even require access to the individual parameter weights or other similar details. In
118 contrast, model-agnostic approaches are designed to operate on any model. They only consider the

119 inputs to a model and corresponding outputs. Another important dimension along which to classify
120 interpretability methods is the scale of interpretation. Global explanations give a high-level overview
121 of how the model performs over its entire mathematical domain. Generally speaking, global
122 explanations sacrifice detailed local precision in favor of global perspective. Local explanations instead
123 are generated for specific examples. They are meant only to be applicable to single predictions of a
124 model and not to offer insights for how the model might classify other examples. The type or modality
125 of data that interpretability methods can be used with, is another crucial factor to consider. *Modality*
126 refers to a particular way in which information is encoded to be presented to a perception/learning
127 system. The following are some examples of different modalities of data: RGB image, infrared image,
128 audio, video, text corpus, GPS trajectory, and time-series sensor data. Typically, models operate on a
129 single modality of data, but they can also be designed to incorporate many modalities at the same time.
130 Lastly, methods for interpretability can be classified based on whether interpretability is achieved by
131 restricting the complexity of the model (intrinsic) by creating a “white-box” or by applying methods
132 that analyze the model after training (post-hoc).

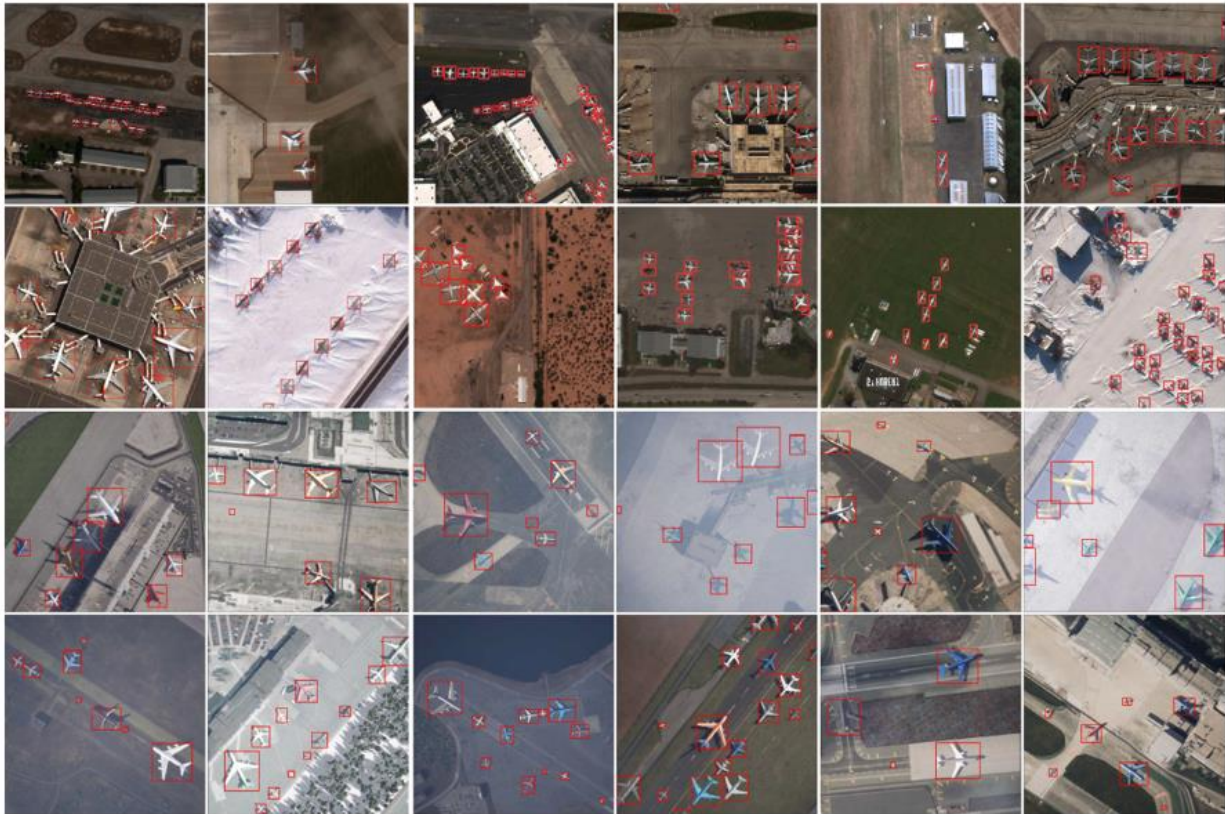
133 In the ATR use case that will be discussed in Section 1.2, local, post-hoc explanation of specific
134 instances is needed to help operators verify that the identified target is, in fact, the correct one.
135 Simplistic versions of ATR utilize only a single modality of data (typically RGB images) but being
136 able to utilize more modalities of information is critical to improving ATR accuracy. In such a system
137 it is useful to be able to use a model-agnostic explanation algorithm since state-of-the-art algorithms
138 will be in a constant state of change and could even be proprietary. A model-agnostic approach to
139 explanations would ensure the explanation system could remain separate and facilitate implementation
140 in such an architecture and help reduce overall maintenance.

141 **1.2 Use Case for Technology Development**

142 This research is driven by the Automatic Target Recognition use case that involves combining the
143 capability of an ATR system with the expertise of human operators; the main goal is to improve the
144 overall target recognition performance by using the strengths of both the human operator and AI/ML
145 algorithms.

146 A typical ATR use case can include a variety of sensors from different domains (space, air, and ground
147 sensors) collecting multi-modal data. Sensor data, for example, can include RGB or Electro Optical
148 imagery. Because of the sheer amount of incoming data, sensor feeds are processed by AI algorithms
149 to highlight potential targets. AI algorithms could potentially come from different vendors with varying
150 degrees of accuracy and reliability. Human analysts (e.g., signal and /or geospatial analysts) then
151 analyze the targeting data generated by the AI systems assessing whether the evidence supports or
152 contradicts their hypothesis and make the target nominations.

153 While performing this work, we consider a scenario in which the ATR system must process imagery
154 data to automatically detect the presence of aircrafts. To this end, we leverage an opensource dataset,
155 called Rareplanes, consisting of 253 Maxar WorldView-3 satellite scenes spanning 112 locations and
156 2,142 km² with 14,700 hand-annotated aircraft data [9]. Figure 1 is an example of some images from
157 the Rareplanes dataset with aircrafts identified by red boxes (i.e., bounding boxes). The dataset can be
158 used to create models that predict the number of engines an aircraft has, the type of wings, the number
159 of tail fins, or the role (civil/military transport, fighter, et cetera), among other things. A pre-trained
160 ML model for detecting aircraft role is also included with the dataset [9]. The Rareplanes dataset and
161 utilities are useful for investigating black box explanations.



162

163 **Figure 1** -- Example images from the Rareplanes dataset. The top 2 rows show examples of real satellite images, the
 164 bottom 2 rows are synthetic images [9].

165 2 Interpretability Methods to Explain Black Box Models

166 There are many survey papers on the broad topic of explainable AI [10],[11], and the more specific
 167 topic of explanation of black box models [5],[12],[13],[14],[15]. This section provides an overview of
 168 selected interpretability methods available in the literature that are model-agnostic and post-hoc. Here
 169 “model-agnostic” indicates that we want to be able to apply these methods to any generic model that
 170 is encountered (i.e., black box model). And “post-hoc” indicates that we are seeking an explanation of
 171 a prediction that has been made by the black box model. Our current focus is on interpretability
 172 methods acting on imagery data as required by the use case under current consideration (Section 1.2).
 173 However, some of the techniques we describe in the following sections might be suitable for other data
 174 modalities. In Sections 2.1 and 2.2 we provide an overview of methods that yield to different types of
 175 explanations including among others saliency maps, that are the most widespread visual explanation
 176 methods, and concepts attribution.

177 2.1 Methods based on Saliency Map

178 Before describing some of the interpretability techniques that can generate post-hoc explanations for
 179 black box models based on imagery data, we provide a definition of saliency map. Given an image that
 180 we want to explain, a saliency map is an image in which a pixel's brightness represents how
 181 salient/important the pixel is [15]. LIME (Local Interpretable Model-agnostic Explanations) [16] and
 182 SHAP (Shapley Additive Explanation) [17] are model-agnostic approaches that provide explanations
 183 of a model prediction by ranking the importance of individual features to a specific prediction,
 184 although, as it will be explained below, they accomplish this in slightly different ways. For imagery

185 data, feature importance can be represented through a saliency map that highlights the contribution of
186 each pixel at the prediction.

187 LIME creates a local, linear, surrogate model around the example for which an explanation is being
188 created. The linear model is created by feeding perturbed versions of the example to the black box
189 model and observing the outputs. For images, the perturbation of the input is achieved by dividing the
190 examples into super-pixel regions, which are groups of neighboring pixels with similar color and
191 brightness. Then a collection of synthetic images is created by replacing random super-pixels of the
192 original image with a uniform, possibly neutral, color. The coefficients of the linear model reflect the
193 importance of each feature to the ultimate prediction.

194 SHAP decomposes the prediction for a given example into the contribution from each individual
195 feature. To estimate the contribution of each feature, perturbed versions of the example are fed to the
196 black box model, and the corresponding outputs are observed. This decomposition is not based on
197 linearizing, and so the parameter rankings are consistent with the original model.

198 In this section we also give an overview of GRADCAM (Gradient-weighted Class Activation Map)
199 that is a post-hoc local explainer for image data [18] that differently from LIME and SHAP is model
200 specific. As it will be explained in Section 3.1 to address the challenge of aligning semantics of the
201 machine generated explanation with human interpretation, we describe the creation of a surrogate
202 white-box approach as a stand-in. Towards that end we built a simple prototype that uses the inputs
203 and outputs of a black box model to train a second explainer model. This second model would serve
204 as a surrogate for the black box model, but since it is *not* black box other algorithms can be used to
205 produce explanations. GRADCAM, based on gradient information of a Convolutional Neural Network
206 (CNN), assigns each neuron a saliency value for the decision of interest and backpropagates this
207 information to the last convolutional layer of the CNN.

208 Approaches like LIME and SHAP generate low-level explanations that in case of imagery inputs are
209 expressed as the presence/absence of pixels that were most important to make a particular decision.
210 However, for AI-aided decision-making systems, the human-in-the-loop that must evaluate the AI
211 decision, usually a domain-expert without technical knowledge in AI, prefers high-level and concept-
212 based explanations that can easily understand and reason with. Such high-level and concept-based
213 explanations are more familiar to humans and more aligned with the human' internal representation of
214 the decision-making problem. Section 2.2 will provide an overview of relevant methods that generate
215 concept-based explanations.

216 **2.2 Methods based on Concept Attribution**

217 In this section we review some of the most relevant methods that generate explanations based on
218 concept attribution. Here concept attribution refers to the ability of a method to quantify how much a
219 concept has contributed to a particular class prediction. Concept-based explanations have the stated
220 goal of using concepts that are “meaningful”, “coherent”, and have “importance” [20]. Here
221 meaningful is defined as having intrinsic meaning to a human; coherent indicates that examples of
222 concepts should be similar to each other and distinct from other concepts; importance signifies that the
223 presence of a concept is necessary to true prediction from examples in a certain class.

224 TCAV (Testing with Concept Activation Vector) [19] quantifies the degree to which hand-selected
225 concepts are important to a classification result. Every concept is represented by a Concept Activation
226 Vector (CAV) that is built by interpreting post-hoc a neural network's internal state in terms of such
227 hand-selected concepts. However, approaches like TCAV require a human to provide hand-labeled

228 examples of concepts. This might introduce human bias in the explanation process by failing to choose
229 the right concepts. Moreover, methods like TCAV that can only leverage concepts that are already
230 labeled and identified by humans, have limited power in discovering other relevant attributes [20].

231 ACE (Automated Concept-based Explanation) [20] is an evolution of TCAV that instead of human-
232 labeled concept data, can automatically discover concepts by breaking up images into segments (e.g.,
233 super-pixels) and by seeing which ones are clustered in the representation space. Then, like in TCAV,
234 ACE quantifies how much these clusters contributed to the prediction of a class.

235 Similarly to ACE, ConceptSHAP [21] targets having concepts consistently clustered to certain
236 coherent spatial regions. For each discovered concept an importance score is defined from a set of
237 Concept vectors (Cs) by utilizing Shapley values.

238 StyleEx [22] introduces a method for automatically discovering visual interpretable attributes and use
239 them for counterfactual explanations. Counterfactual explanations provide alternative inputs, where a
240 small set of attributes is altered, and the different classification outcomes are observed. For each
241 discovered attribute, StyleEx generates a counterfactual example showing how manipulating such
242 attribute affects the classifier prediction (*Had the input x been $\neg x$ then the classifier output would have*
243 *been $\neg y$ instead of y*). To generate visual counterfactual explanations StyleEx leverages a StyleGAN
244 [23] trained to discover classifier-related attributes.

245 The concept-based explanation literature presented in this section is lacking in some respects. On one
246 hand, techniques like TCAV require resources to generate human hand-labeled examples of concepts
247 and its discovery power is rather limited. On the other hand, techniques like ACE and StyleEx can
248 automatically discover relevant concepts but there is no guarantee that such concepts will be human
249 interpretable. Research related to such methods have only found evidence that the discovered concepts
250 *tend* to be meaningful to people.

251 Relative to the literature described in Sections 2.1 and 2.2, there are still many research challenges that
252 need to be addressed in pursuit of creating user-intuitive and effective explanations of black box
253 machine learning models. Section 3 discusses key areas requiring further investigation for improving
254 the quality of model-agnostic explanations.

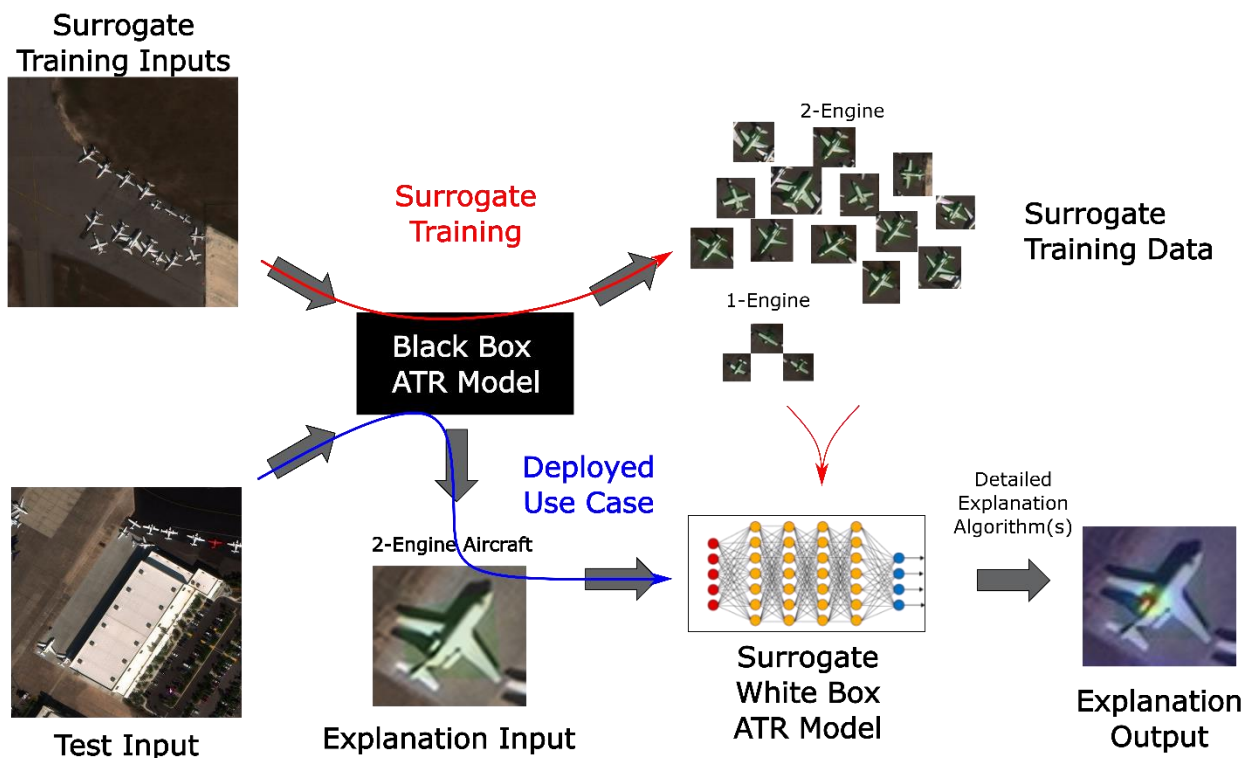
255 **3 Toward User-Intuitive Explanations**

256 Treating the model as a black box while producing explanations is convenient because it allows the
257 predictive and explanation systems to be entirely decoupled. However, most black box explanation
258 methods are currently reliant on the inherent interpretability of the features that they use. That is, they
259 assume that the semantics of the explanations are aligned with human interpretations. The ease of
260 interpreting the explanation is directly related to the extent to which the features are interpretable. This
261 may work generally, but in performance-critical applications a more intentional selection of
262 semantically meaningful features is desirable, perhaps necessary. Such an approach would help ensure
263 that produced explanations would be more meaningful to an operator. To address this challenge of
264 aligning semantics of the machine generated explanation with human interpretation, we first describe
265 the creation of a surrogate white-box approach as a stand-in (Section 3.1), and subsequently describe
266 the concept of a semantic alignment method (Section 3.2) based on a neuro-symbolic approach.

267 **3.1 Creating a Surrogate White Box Model**

268 Methods like LIME and SHAP reviewed in Section 2.1 are limited to ranking feature importance
 269 (although “features” can be quite a flexible concept). One way to address this is to use the inputs and
 270 outputs of the black box model to train a second explainer model. This second model would serve as a
 271 surrogate for the black box model, but since it is *not* black box other algorithms can be used to produce
 272 explanations.

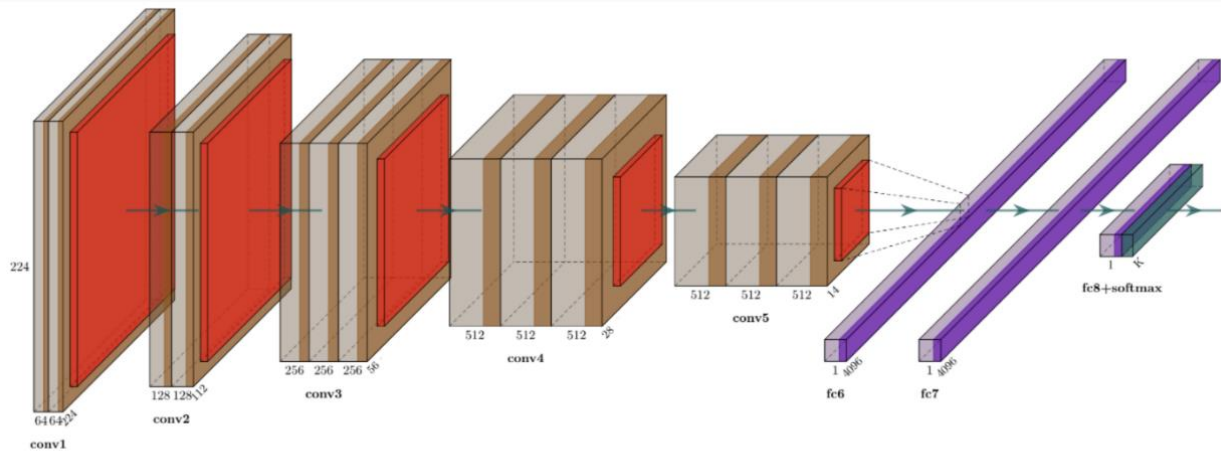
273 A simple prototype of this system was created to investigate the possibilities. The main goal was to
 274 create a white-box model based on outputs from the black box model. This was done in the following
 275 way: first, we used labeled training data from the Rareplanes dataset; specifically, images of aircraft
 276 along with their classification of how many engines they have (an integer 0 through 4). This set of
 277 images was then used to train a convolutional neural network (CNN) to use those cropped images and
 278 classify the number of engines. Finally, a white-box, model-aware, explanation method known as
 279 Grad-CAM [18] was used to produce saliency maps illustrating the regions of images that were most
 280 influential for making a classification. This process is depicted in Figure 2.



281

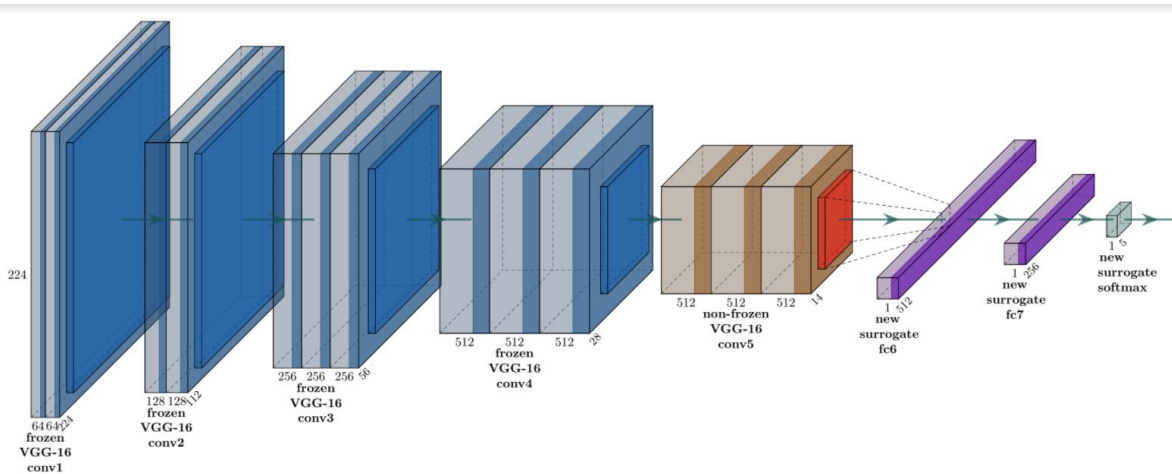
282 **Figure 2** -- Diagram illustrating the training process for the surrogate model as well as an example of how the model
 283 could be used after being deployed.

284 The surrogate white box model used for this prototype was based on the VGG-16 model originally
 285 designed for the ImageNet dataset [24]. The architecture of the original VGG-16 model is shown in
 286 Figure 3; for our prototype specialized model a modified version of the VGG-16 model was created
 287 (Figure 4). In essence convolutional layers 1-4 were “frozen” (kept their weights from being trained
 288 on ImageNet). The 5th convolutional layer had its weights re-trained on the surrogate training data set
 289 discussed in Figure 2. The final fully connected layers had their sizes changed to better fit the specific
 290 classification problem. Note that this description highlights what we did for the prototype system that
 291 we built and is not meant to be generally prescriptive for future work; it simply demonstrates that it is
 292 possible to build a surrogate white box model as a means to utilize different non-black box explanation
 293 generation algorithms.



294
295

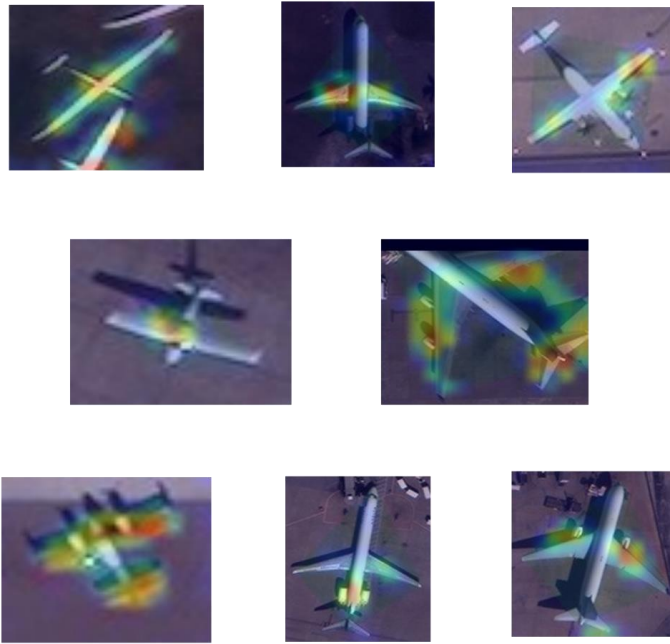
Figure 3 -- Diagram showing the layout of the VGG-16 architecture [24].



296

297 **Figure 4** -- Surrogate white box model reuses the architecture from the original VGG-16 model. Layers ‘conv1’, ‘conv2’,
298 ‘conv3’, and ‘conv4’, were left unmodified (i.e., “frozen”). The ‘conv5’ layer kept the same structure, but the weights were
299 modified after training on the labeled images discussed. Layers, ‘fc6’ and ‘fc7’ had their size changed; ‘fc8’ was removed,
300 and the final ‘softmax’ layer was changed to accommodate the 0-4 engine classification output of the Rareplanes dataset.

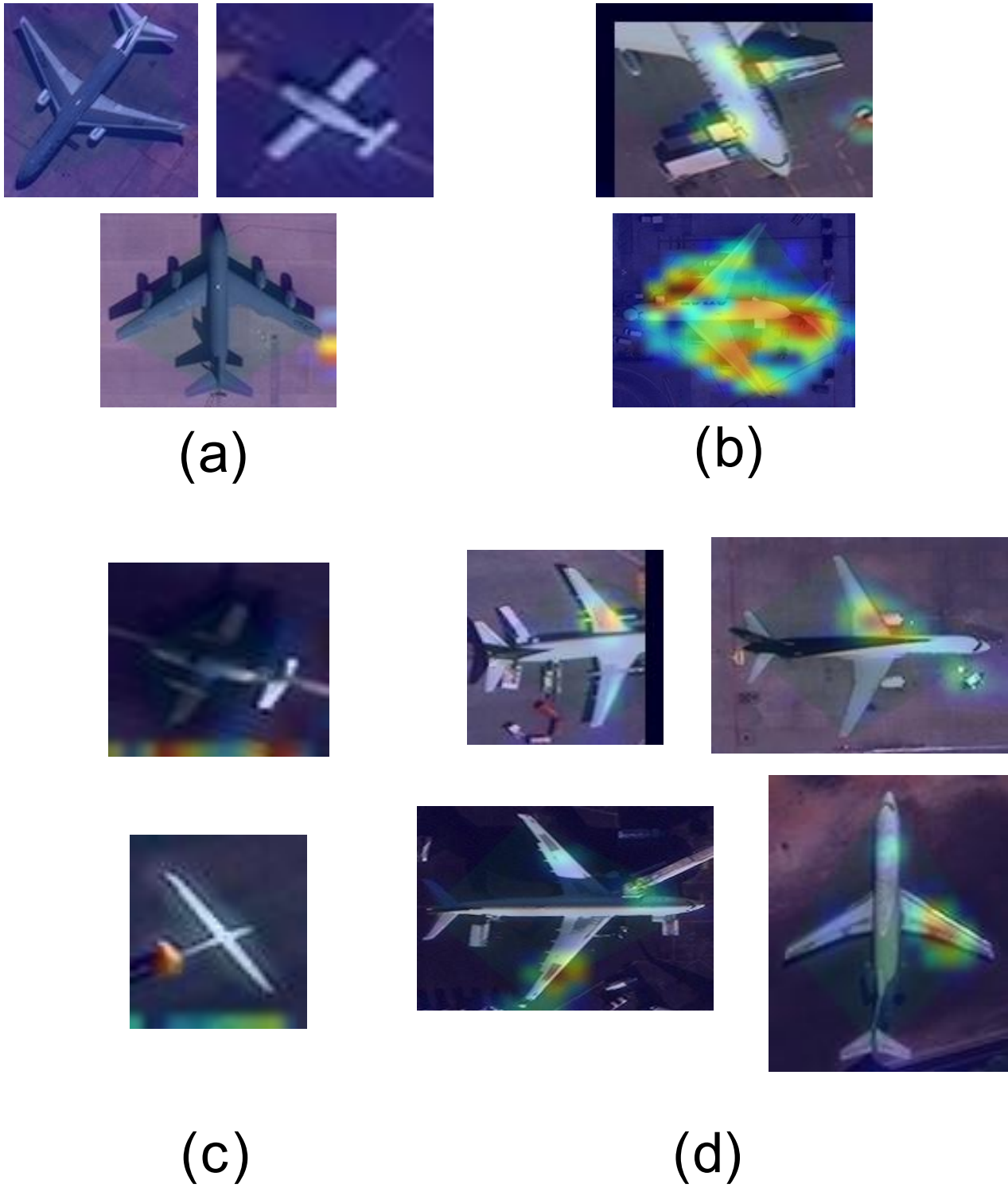
301 We then applied the Grad-CAM explanation method to produce saliency maps illustrating the regions
302 of images that were most influential for making a classification. It is important to note that while we
303 used Grad-CAM in this prototype, it is also possible to apply a host of other explanation methods. This
304 serves as a proof of concept that it’s possible to train a surrogate white box model that enables the
305 application of non black box explanation algorithms to explain the decisions of a black box model.
306 Figure 5, and Figure 6 contain example explanations of a surrogate white-box model trained to mimic
307 the Rareplanes classifications of the number of engines as illustrated in Figure 2. Examples in Figure
308 5 generally demonstrate promising behavior, those in Figure 6 highlight some of the existing
309 limitations. Figure 7, and Figure 8 contain results from similar surrogate models created to explain the
310 role of the aircraft, and to classify the number of tail fins respectively.



311

312
313

Figure 5 -- Selected Grad-CAM results from the surrogate white box model. The images uniformly indicate that the important areas for classifying the number of engines are the wings, tail, and engines of a plane.

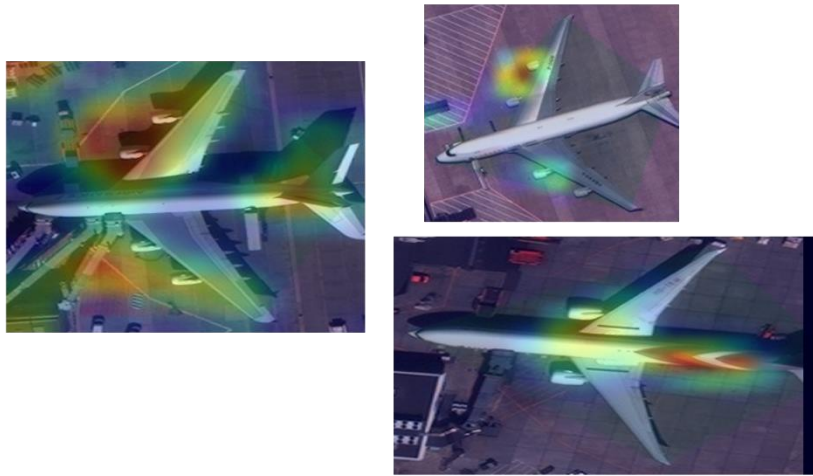


314

315 **Figure 6** -- Selected results from Grad-CAM analysis of the surrogate model. These images represent unexpected results
 316 sometimes given by the algorithm. In (a) we see that sometimes the gradient is not very informative. The images in (b)
 317 highlight that sometimes the algorithm focused on too many features, or seemingly irrelevant features. Sometimes the
 318 gradient was at the edge of an image with almost not overlap on the plane (c). Finally, at other times the model seems to
 319 focus primarily on a single wing (d).

320

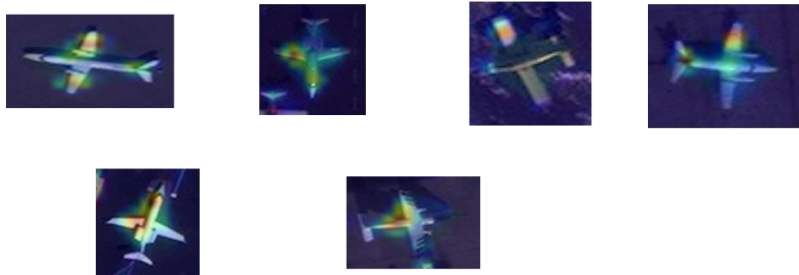
Large Transport



Medium Transport



Small Transport



321

322
323

Figure 7 – Example explanations for prediction of civil role. The Rareplanes dataset has very few military aircraft (and thus has poor performance) so we don't display those results.

324



Figure 8 – Example explanations for a surrogate white-box classifier that predicts the number of tail fins. This surrogate model doesn’t seem to be working correctly. Notice that there are few images that have hot-spots where GradCAM shows a high gradient.

3.2 Semantically Aligning the Explanation Model: A Neuro-symbolic Approach

Current machine learning methods can provide high classification accuracy but low human interpretability because the intermediary “hidden” layers are abstract statistical transformations. Indeed, in ML-aided decision-making tasks, the human-in-the-loop, usually a domain-expert without technical ML knowledge, prefers high-level concept-based explanations instead of low-level explanations based on model features. In order to create meaningful explanations for a black box model, that can align with user’ internal representation of the problem, concept-based explanations can be leveraged. However, as discussed in Section 2.2 current methods based on concept attribution either require resources to generate hand-labeled examples of concepts or, in case such attributes can be automatically discovered, there is no guarantee that such concepts will be human interpretable. Most recently neuro-symbolic AI has been proposed as tool for achieving more explainable, transparent, and trustworthy systems. The neuro-symbolic approach is mostly motivated by the complementary strengths of sub-symbolic computations (i.e., data-driven AI) and symbolic representations (i.e., knowledge-driven AI) that could lead to designing hybrid, more intelligent AI systems [25][24]. More specifically, the neuro-symbolic approach offers more explainable, transparent and trustworthy systems by complementing sub-symbolic approaches and their ability to deal with large amounts of data, to handle noise, and to capture the richness of perceptual data, with symbolic methods, allowing to encode knowledge in the form of language-like, structured propositions that can be endlessly recombined to allow high-level reasoning across tasks and domains [26].

348 Knowledge graphs have been proposed as tools to support the explanation of machine learning models
349 [26],[27],[28]. In the context of image classification tasks, early work focused on manually creating an
350 ontology capturing spatial concepts, colors, textures, and their relationships, and incorporating it in an
351 object recognition classifier [29]. The domain knowledge captured in the ontology was leveraged to
352 foster model transparency acting as a user-friendly intermediate between the classifier and the end-user
353 [29]. Most recently, existing large-scale and open-source knowledge graphs, such as OpenCyc¹,
354 ConceptNet² and DBpedia³, have been embedded in machine learning models to increase model
355 transparency [30],[31],[32]. For example, background knowledge from ConceptNet was leveraged to
356 explain objects in images with associated sentences in the form of captions [33]. A sentence-based
357 image retrieval problem was leveraged to demonstrate that keywords in the captions that did not have
358 a visual detector available, could be explained by leveraging the concepts and relations in the
359 knowledge graph connected to such keywords.

360 Most of the neuro-symbolic approaches proposed so far do require a manual step of extracting
361 knowledge sources. Indeed, one open challenge in knowledge-based explainable systems is the ability
362 to automatically extract knowledge from graphs. In current approaches the identification of the correct
363 portion of information in the graph to generate explanations requires using a human expert.

364 **4 Metrics for Explanation Evaluation**

365 This section describes relevant metrics that can be leveraged for measuring the effectiveness of the
366 explanations generated by the techniques described in the previous sections of this report. For our
367 purposes, the effectiveness of explanations must be assessed by evaluating how helpful the explanation
368 is to the human end user. Such assessment requires human-in-the-loop experiments based on a variety
369 of metrics and fall under the category of Human-grounded Metrics [34].

370 The DARPA XAI program formulated a model of explanation inspired by psychological principles
371 [35]. The model describes the process underlying the interaction between a human end user and an
372 explainable AI system as an opportunity to introduce metrics for assessing the human + AI system
373 performance. Following the DARPA XAI model, initial instructions on how to use an AI system enable
374 a human end-user to form an early mental model of the task and the AI system. As the end-user starts
375 using the AI system, explanations are generated to help the human refine his/her mental model. This
376 in turn should take to better performance of the human + AI system, and to the establishment of
377 appropriate trust and reliance in the AI system. To evaluate the performance of the human + AI system,
378 five measurement categories were proposed in [35].

379 The “*Goodness Criteria*” provides a means to assess the goodness of explanations, based on factors
380 revolving around clarity and precision. A “goodness checklist” that can be used to either design
381 “goodness” into the explanations, or to evaluate the a-priori goodness of the explanations that the
382 system generates is described.

383 The “*Test of Satisfaction*” enables the assessment of whether users are satisfied by the received
384 explanations. This is an a-posteriori judgment of explanations that leverages terms such as
385 understandability, feeling of satisfaction, sufficiency of detail, completeness, usefulness, accuracy, and
386 trustworthiness, included into a Likert scale for review and evaluation.

¹ <http://www.cyc.com/opencyc/a>

² <http://conceptnet.io/>

³ <https://www.dbpedia.org>

387 The “*Test of Understanding*” gives information on how well users understand the AI systems (i.e.,
388 mental model). Methods to elicit human users’ mental models, resulting in data that can be easily
389 scored, categorized, or analyzed are described. The intent of such methods is to provide some sort of
390 structure or “scaffolding” that supports the user in explaining their reasoning process while solving a
391 task. Examples of eliciting users’ mental models include a “*Retrospection Task*” for which probe
392 questions are presented to participants about their reasoning just after the reasoning task has been
393 performed (e.g., “*Can you describe the major components or steps in the [software system,
394 algorithm]?*”). Another example of eliciting users’ mental models is a “*Prediction Task*” during which
395 users are presented with test-cases, and they are asked to predict the results and then explain why they
396 thought they would obtain those results (e.g., “*What will the [software system, algorithm] do next?*”,
397 “*How do I intervene?*”). In a “*Self-explanation Task*” users express their own understanding by
398 generating self-explanations. This helps learners to refine their knowledge. Last, a “*Think-aloud
399 Solving Task*” participants think aloud while they solve a task.

400 The “*Test of Performance*” provides information on how the human-AI system perform with the intent
401 to understand if there is an improvement in the user’s decision and task performance. According to this
402 model, the user performance will improve as a result of receiving satisfying explanations, and it will
403 be a function of the quality of his/her mental model.
404

405 The last measurement proposed in the DARPA XAI program is about “*Appropriate Trust and
406 Reliance*” providing information on whether the user’s trust and reliance on the AI system are
407 appropriate. Currently, the assessment of trust in such systems is predominantly done via human self-
408 reported behavior. However, rather than looking at scales for measuring interpersonal trust, the interest
409 here was in scales designed to assess human trust in AI systems. Typical scales include questions to
410 assess trust (e.g., “*Do you trust the machine’s outputs?*”) and reliance (e.g., “*Would you follow the
411 machine advice?*”). Other available scales are highly specific to the context and application of interest
412 (e.g., assistive robotic technology for elderly). A trust scale that incorporates items from other relevant
413 scales and that might be used in the XAI context has been defined in [35]. According to the outcome
414 of the XAI Program, trust assessment should be a repeat measure that requires multiple measures taken
415 over time and integrated for overall evaluations of human user + system performance.

416 **4.1 Future Research Needs**

417 For the research described in Section 3, new effective evaluation metrics based on trustworthiness and
418 acceptance will need to be defined. Future research directions in this area can incorporate technologies
419 for continuous monitoring and automatic assessment of human trust in AI systems. Indeed, trust
420 depends on time-varying factors that can influence the human decision-making process during
421 interactions with AI systems [36]; the intended user should learn under what conditions the system
422 fails and it succeeds in accomplishing the user’s goals [7]. Further, trust may also be built or calibrated
423 “on the job” as operation time increases, facilitated by the incorporated continuous monitoring and
424 automatic assessment technologies. Interesting directions to explore include leveraging
425 psychophysiological measurements to collect data that, when used as input to classification algorithms,
426 map continuous data to a categorical trust level.

427 **5 Conclusion**

428 This work has focused on methods for making AI/ML models, leveraged by Automatic Target
429 Recognition systems for example, more interpretable to human operators by employing AI explainer
430 generation technologies. In ATR it is necessary for the human operator to be able to verify that

431 identified targets meet the necessary criteria; user-intuitive and accurate explanations will help them
432 to do so. Model-agnostic and post-hoc algorithms were identified as promising technologies. There are
433 still several technical challenges that exist before this technology can be considered capable of meeting
434 all desired characteristics. Approaches like LIME and SHAP generate low-level explanations that in
435 case of imagery inputs are expressed as the presence/absence of pixels that were most important to
436 make a particular decision. However, for AI-aided decision-making systems, the human-in-the-loop
437 that must evaluate the AI decision, usually a domain-expert without technical knowledge in AI, prefers
438 high-level and concept-based explanations that can easily understand and reason with. Such high-level
439 and concept-based explanations are more familiar to humans and more aligned with the human' internal
440 representation of the decision-making problem. The concept-based literature has made significant
441 progress in methods that can automatically discover relevant concepts and quantify how much such
442 concepts contribute to a particular class prediction. However, there are still challenges that need to be
443 addressed. On one hand, techniques like TCAV require resources to generate human hand-labeled
444 examples of concepts and its discovery power is rather limited. Techniques like ACE and StyleX can
445 automatically discover relevant concepts but there is no guarantee that such concepts will be human
446 interpretable. Research related to such methods have only found evidence that the discovered concepts
447 *tend* to be meaningful to people. To address this challenge of aligning semantics of the machine
448 generated explanation with human interpretation, we describe the concept of a semantic alignment
449 method and approaches that can be leveraged for that. Recently neuro-symbolic AI has been proposed
450 as tool for achieving more explainable, transparent, and trustworthy systems. The neuro-symbolic
451 approach is mostly motivated by the complementary strengths and weaknesses of sub-symbolic
452 computations (i.e., data-driven AI) and symbolic representations (i.e., knowledge-driven AI) that could
453 lead to designing hybrid, more intelligent AI systems. We also describe relevant metrics that can be
454 leveraged for measuring the effectiveness of the explanations generated by the techniques described in
455 this work. For our purposes, the effectiveness of explanations must be assessed by evaluating how
456 helpful the explanation is to the human end user. Such assessment requires human-in-the-loop
457 experiments based on a variety of metrics and fall under the category of Human-grounded Metrics. For
458 the research described in this work, new effective evaluation metrics based on trustworthiness and
459 acceptance will need to be defined. Future research directions in this area can incorporate technologies
460 for continuous monitoring and automatic assessment of human trust in AI systems.

461 **6 Acknowledgments**

462 This project was funded by Raytheon Intelligence & Space Internal Research and Development
463 (POCs Frank Tanner and Laura Strater).

464 **7 References (APA formatting)**

- 465 [1] Bocskor, A., Hunyadi, M., & Vince, D. (2017). National Academies of Sciences, Engineering, and Medicine (2015)
466 The Integration of Immigrants into American Society. Washington, DC: The National Academies Press. 458 pages.
467 INTERSECTIONS: EAST EUROPEAN JOURNAL OF SOCIETY AND POLITICS, 3(3), 157-161.
- 468 [2] Endsley, M. R., Bolté, B., & Jones, D. G. (2003). Designing for situation awareness: An approach to user-centered
469 design. CRC press.
- 470 [3] Chen, J. Y., Procci, K., Boyce, M., Wright, J., Garcia, A., & Barnes, M. (2014). Situation awareness-based agent
471 transparency. Army research lab aberdeen proving ground md human research and engineering directorate.
- 472 [4] Gunning, D., Vorm, E., Wang, J. Y., & Turek, M. (2021). DARPA's explainable AI (XAI) program: A retrospective.
473 Applied AI Letters, 2(4), e61.
- 474 [5] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). A survey of methods for
475 explaining black box models. ACM computing surveys (CSUR), 51(5), 1-42.
- 476 [6] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint
477 arXiv:1702.08608.

- 478 [7] Board, D. I. (2019). AI Principles: Recommendations on the Ethical Use of Artificial Intelligence by the Department
479 of Defense: Supporting Document. [https://media.defense.gov/2019/Oct/31/2002204458/-1/-](https://media.defense.gov/2019/Oct/31/2002204458/-1/-1/0/DIB_AI_PRINCIPLES_PRIMARY_DOCUMENT.PDF)
480 [1/0/DIB_AI_PRINCIPLES_PRIMARY_DOCUMENT.PDF](https://media.defense.gov/2019/Oct/31/2002204458/-1/-1/0/DIB_AI_PRINCIPLES_PRIMARY_DOCUMENT.PDF)
- 481 [8] Linardatos, P., Papastefanopoulos, V., & Kotsiantis, S. (2021). Explainable ai: A review of machine learning
482 interpretability methods. *Entropy*, 23(1), 18.
- 483 [9] Rareplanes: Synthetic data takes flight. In *Proceedings of the IEEE/CVF Winter Conference on Applications of*
484 *Computer Vision* (pp. 207-217).
- 485 [10] Burkart, N., & Huber, M. F. (2021). A survey on the explainability of supervised machine learning. *Journal of Artificial*
486 *Intelligence Research*, 70, 245-317.
- 487 [11] Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M., & Kagal, L. (2018, October). Explaining explanations: An
488 overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and*
489 *advanced analytics (DSAA)* (pp. 80-89). IEEE.
- 490 [12] Mittelstadt, B., Russell, C., & Wachter, S. (2019, January). Explaining explanations in AI. In *Proceedings of the*
491 *conference on fairness, accountability, and transparency* (pp. 279-288).
- 492 [13] Joshi, G., Walambe, R., & Kotecha, K. (2021). A review on explainability in multimodal deep neural nets. *IEEE*
493 *Access*.
- 494 [14] Vermeire, T., Laugel, T., Renard, X., Martens, D., & Detyniecki, M. (2021). How to choose an Explainability Method?
495 Towards a Methodical Implementation of XAI in Practice. *arXiv preprint arXiv:2107.04427*.
- 496 [15] Bodria, F., Giannotti, F., Guidotti, R., Naretto, F., Pedreschi, D., & Rinzivillo, S. (2021). Benchmarking and survey of
497 explanation methods for black box models. *arXiv preprint arXiv:2102.13076*.
- 498 [16] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). " Why should i trust you?" Explaining the predictions of any
499 classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data*
500 *mining* (pp. 1135-1144).
- 501 [17] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in neural*
502 *information processing systems*, 30.
- 503 [18] Selvaraju, R. R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., & Batra, D. (2016). Grad-CAM: Why did you say
504 that? *arXiv preprint arXiv:1611.07450*.
- 505 [19] Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., & Viegas, F. (2018, July). Interpretability beyond feature
506 attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*
507 (pp. 2668-2677). PMLR.
- 508 [20] Ghorbani, A., Wexler, J., Zou, J. Y., & Kim, B. (2019). Towards automatic concept-based explanations. *Advances in*
509 *Neural Information Processing Systems*, 32.
- 510 [21] Yeh, C. K., Kim, B., Arik, S., Li, C. L., Pfister, T., & Ravikumar, P. (2020). On completeness-aware concept-based
511 explanations in deep neural networks. *Advances in Neural Information Processing Systems*, 33, 20554-20565.
- 512 [22] Lang, O., Gandelsman, Y., Yarom, M., Wald, Y., Elidan, G., Hassidim, A., ... & Mosseri, I. (2021). Explaining in
513 Style: Training a GAN to explain a classifier in StyleSpace. In *Proceedings of the IEEE/CVF International Conference*
514 *on Computer Vision* (pp. 693-702).
- 515 [23] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., & Aila, T. (2020). Analyzing and improving the image
516 quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 8110-
517 8119).
- 518 [24] Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv*
519 *preprint arXiv:1409.1556*.
- 520 [25] Van Harmelen, F., & Teije, A. T. (2019). A boxology of design patterns for hybrid learning and reasoning systems.
521 *arXiv preprint arXiv:1905.12389*.
- 522 [26] Tiddi, I., & Schlobach, S. (2022). Knowledge graphs as tools for explainable machine learning: A survey. *Artificial*
523 *Intelligence*, 302, 103627.
- 524 [27] Lecue, F. (2020). On the role of knowledge graphs in explainable AI. *Semantic Web*, 11(1), 41-51.
- 525 [28] Futia, G., & Vetrò, A. (2020). On the integration of knowledge graphs into deep learning models for a more
526 comprehensible AI—Three Challenges for future research. *Information*, 11(2), 122.
- 527 [29] Mailliot, N. E., & Thonnat, M. (2008). Ontology based complex object recognition. *Image and Vision Computing*,
528 26(1), 102-113.
- 529 [30] Sarker, M. K., Xie, N., Doran, D., Raymer, M., & Hitzler, P. (2017). Explaining trained neural networks with semantic
530 web technologies: First steps. *arXiv preprint arXiv:1710.04324*.
- 531 [31] Daniels, Z. A., Frank, L. D., Menart, C. J., Raymer, M., & Hitzler, P. (2020, April). A framework for explainable deep
532 neural models using external knowledge graphs. In *Artificial Intelligence and Machine Learning for Multi-Domain*
533 *Operations Applications II* (Vol. 11413, p. 114131C). International Society for Optics and Photonics.
- 534 [32] Marino, K., Salakhutdinov, R., & Gupta, A. (2016). The more you know: Using knowledge graphs for image
535 classification. *arXiv preprint arXiv:1612.04844*.

- 536 [33]Icarte, R. T., Baier, J. A., Ruz, C., & Soto, A. (2017). How a general-purpose commonsense ontology can improve
537 performance of learning-based image retrieval. arXiv preprint arXiv:1705.08844.
- 538 [34]Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint
539 arXiv:1702.08608.
- 540 [35]Hoffman, R. R., Mueller, S. T., Klein, G., & Litman, J. (2018). Metrics for explainable AI: Challenges and prospects.
541 arXiv preprint arXiv:1812.04608.
- 542 [36]Lee, J. D., & See, K. A. (2004). Trust in automation: Designing for appropriate reliance. *Human factors*, 46(1), 50-80.
- 543 [37]Brett W. Israelsen and Nisar R. Ahmed. 2019. "Dave...I Can Assure You...That It's Going to Be All Right..." A
544 Definition, Case for, and Survey of Algorithmic Assurances in Human-Autonomy Trust Relationships. *ACM Comput.*
545 *Surv.* 51, 6 (January 2019), 113:1–113:37.